

# **Data Splash!**

## **An Educational Game about Machine Learning**

Andrew Gray

445348

Submitted to Swansea University in fulfilment  
of the requirements for the Degree of Master of Science



**Swansea University**  
**Prifysgol Abertawe**

Department of Computer Science  
Swansea University

September 20<sup>th</sup>, 2020

# Declaration

This work has not been previously accepted in substance for any degree and is not being concurrently submitted in candidature for any degree.

Signed ..... (candidate)

Date .....

# Statement 1

This thesis is the result of my own investigations, except where otherwise stated. Other sources are acknowledged by footnotes giving explicit references. A bibliography is appended.

Signed ..... (candidate)

Date .....

# Statement 2

I hereby give my consent for my thesis, if accepted, to be available for photocopying and for inter-library loan, and for the title and summary to be made available to outside organisations.

Signed ..... (candidate)

Date .....



*I would like to dedicate this work to the Hypnotoad.  
All glory to the Hypnotoad.*



# Abstract

In your abstract you should aim to summarize the core contributions of your work in the context of the problem domain. Start by outlining the domain and the problems posed within it. Discuss how the methods you focus on approach the relevant problems. You should end your abstract by concretely stating the tangible outputs and deliverables you have created in order to complete your work on this document, and whether those outputs represent and improvement or alternative approach to existing methods.

Your abstract should be a couple or so paragraphs long, and roughly approximate the order and flow you then use for structuring the main document. If a viewer has read your abstract then they should already understand at a high level what it is you have created and delivered, and whether it is better than or comparable to existing methods. If your project is driven by a research hypothesis then the reader should know what that is at a high level from this section. Reading on, little should surprise the viewer.

For paper submission of your thesis you should physically sign your name on each of the above declaration statements and date them in black ink. For digital submissions you should sign and date them digitally using a touch or stylus input if available. There are pieces of software that allow you to write directly on PDF documents, or alternatively you can bring a signature into your document as a figure with a transparent or white background. If you do not have a stylus input / tablet like device you should ask your supervisor, as many in the department do their grading / work on digital tablets.





# Acknowledgements

This is an opportunity to acknowledge and thank those who have supported you throughout your studies. Friends and colleagues who you have studied alongside, your families, and your mentors within the department are the usual suspects.



# Contents

<b>List of Tables</b>	<b>vii</b>
<b>List of Figures</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Overview of the Problem . . . . .	1
1.2 Overview of the Solution . . . . .	1
1.3 Motivations . . . . .	1
1.4 Aims . . . . .	2
1.5 Contributions . . . . .	2
1.6 Thesis Overview . . . . .	2
<b>2 Background &amp; Literature Review</b>	<b>3</b>
<b>3 Methodology</b>	<b>5</b>
3.1 Overview of Application . . . . .	5
3.2 Overview of Specific Game Components . . . . .	6
3.3 Evaluation of Application . . . . .	6
<b>4 Implementation</b>	<b>9</b>
4.1 Tools . . . . .	9
4.2 Frameworks . . . . .	13
4.3 Packages . . . . .	14
4.4 IDE . . . . .	14
4.5 Intricacies of the Game Components . . . . .	16
4.6 Example user stories (A UML term for case studies or example playthroughs) . . . . .	20

<b>5</b>	<b>Conclusions and Future Work</b>	<b>21</b>
5.1	Contributions . . . . .	21
5.2	Future Work . . . . .	21
	<b>Bibliography</b>	<b>23</b>
	<b>Appendices</b>	<b>25</b>
A	Implementation of a Relevant Algorithm	27
B	Supplementary Data	29

# List of Tables

# List of Figures

- 4.1 A comparison between Java, Python and C++ to print an output to the console. [1] . 10

# Chapter 1

## Introduction

**As part of our Masters of Science accreditation, we must complete a research thesis. It has been decided upon to create an educational game centred around Machine Learning (ML), due to the authors desire to gain a deeper understanding of ML and their previous experiences as being a secondary school teacher.**

### 1.1 Overview of the Problem

### 1.2 Overview of the Solution

### 1.3 Motivations

Large documents can become cumbersome to work with and format consistently. Sensibly chosen aesthetic cues are important to help imply structure and can greatly aid the reader in understanding your work. The accompanying LaTeX template uses abstraction to hide the formatting from the author during content preparation, allowing for consistent styling to be applied automatically during document compilation. In this Google Docs theme it is the responsibility of the author to manually adhere to the styling laid out in this template.

## 1.4 Aims

### 1.4.1 Objective

In this document we present a tutorial on thesis creation and typesetting, and discuss topics such as literature surveying and proper citation.

## 1.5 Contributions

The main contributions of this work can be seen as follows:

- **A LaTeX thesis template**

Modify this document by adding additional TeX files for your top level content chapters.

- **A typesetting guide of useful primitive elements**

Use the building blocks within this template to typeset each part of your document. Aim to use simple and reusable elements to keep your LaTeX code neat and to make your document consistently styled throughout.

- **A review of how to find and cite external resources**

We review techniques and resources for finding and properly citing resources from the prior academic literature and from online resources.

## 1.6 Thesis Overview

The remainder of chapter 1 outlines the document structure and the key contributions of this work is organized as follows. Chapter ?? reviews techniques for finding and properly citing external resources from the academic literature and online. In chapter ?? we show examples of how to typeset different types of content, such as internal references, figures, code listings, and tables. And lastly in chapter 5 we summarize the main contributions and key points to take away from this template.



## **Chapter 2**

# **Background & Literature Review**



## Chapter 3

# Methodology

The university has subscriptions to a vast number of major academic journals spanning a wide range of subject areas. By accessing the internet from a university network connection (Eduroam or Ethernet), the paywalls of many journals will simply vanish without any need for login credentials.

### 3.1 Overview of Application

When you are working from outside of the university then connecting to an on campus machine via remote desktop (RemoteDesktopProtocol, TeamViewer, ect) or via port forwarding (ssh, ssh tunnel, ect) can allow you to access papers that would otherwise be behind a paywall.

If you do not have individual access to a machine that is exposed for ssh on the university network you can always use the computers in Linux Lab CF204<sup>1</sup> for the purpose of setting up an ssh port tunnel to proxy your internet through. These machines have fixed IPv4 addresses and respond to ssh using your student account credentials. While in use your internet will be routed<sup>2</sup> to the university and then out to the internet, granting you transparent access to journals without a paywall.

---

<sup>1</sup>One caveat of using computer lab machines for remote tunnelling is that a environmentally conscious student who has worked late in the computer lab might choose to switch off the machine you were using...

<sup>2</sup>Painfully slowly.

#### 3.1.1 Design

### 3.2 Overview of Specific Game Components

The internet is big [2]. Knowing how to phrase a question to a search engine is therefore an invaluable skill. If the request is simple enough, even a poorly structured query will likely return usable results. For more difficult to find resources you can leverage the language of the search engine to gather relevant papers and resources for your research more efficiently.

<https://www.gwern.net/Search>

“Internet Search Tips” [3] provides an excellent review of methods and tips for scouring the internet for hard to find resources. You will also be less likely to get caught behind journal paywalls when working remotely without a tunnel as your queries can be made to look for raw pdfs that are often released by the authors directly.

#### 3.2.1 Game Arena

#### 3.2.2 Free Play

#### 3.2.3 Learning Zone

#### 3.2.4 Awards Zone

### 3.3 Evaluation of Application

BibTeX is a language for specifying resource citations. Every time you access and read an academic paper, take code from an online repository, or source the media such as images from existing works you should create a BibTeX entry in a file that you keep throughout your research. Software such as Mendeley [4] can help automate the process of building your BibTeX library of citations.

```
1 @INPROCEEDINGS{kaj86,
2   author   = {Kajiya, James T.},
3   title    = {The Rendering Equation},
4   booktitle = {Proceedings of the 13th Annual Conference on Computer Graphics
5               and Interactive Techniques},
6   year     = {1986},
7   series   = {SIGGRAPH '86},
8   pages    = {143--150},
9   address  = {New York, NY, USA},
```

```
9 | publisher = {ACM},  
10 | isbn      = {0-89791-196-2},  
11 | numpages  = {8},  
12 | acmid     = {15902}  
13 | }
```

Listing 3.1: An example BibTeX entry for an academic paper published in conference proceedings [5].

The BibTeX code listing above (listing 3.1) shows an example of how to cite an academic paper, in this case one of the central papers in Computer Graphics research. The key **kaj86** is an arbitrary name chosen as a meaningful identifier for the resource. In the document text we can call on this resource as an inline citation using the LaTeX command `\cite{kaj86}` which produces [5] at the location it is called. As long as a citation has been used at least once somewhere within the document then a formatted full citation will be created in the bibliography at the end of the document with the same citation number that is shown inline.

It is considerably easier to be disciplined in methodically taking note of the resources you access and make use of as you access them, than it is to try and hunt them all down again at the time you need to write about them in your document. Invest time in being organized and consistent up front and it will be easier when you come to write up.

### 3.3.1 User Study



## Chapter 4

# Implementation

### 4.1 Tools

#### 4.1.1 Programming Languages

For the implementation of our application, three primary programming languages deemed to be best suited for development. Apple's Swift programming language [6] got considered early on, due to the author's familiarisation with the programming language. The programming language gets used for creating applications for Apple's mobile and desktop operating systems, and with 1.5 billion [7] iOS devices in circulation, that was a lot of potential users. Additionally, Apple's iOS devices are prevalent within most educational settings, with Apple's iPad being one of the primary go-to devices. However, due to the language not supporting key frameworks required, or providing similar alternatives, the decision to not use this language got made.

We then got presented with three main options to use, Python, R and HTML, CSS and JavaScript.

Python is a very popular programming language [8, 9], it is fast, easy-to-use, and easy-to-deploy programming language that gets widely used to develop scalable applications. Examples include YouTube, Instagram, Pinterest and SurveyMonkey [10]. The Python Software Foundation state that Python is a high-level, object-orientated (OOP), interpreted language with dynamic semantics. Due to the language being a high-level, it has many built-in data structures. These features, along with the dynamic typing and dynamic binding together make Python attractive to development teams working in a Rapid Application Development (RAD). As Python is an extracted level above the C language [11], Python can get used as the glue that connects existing components, as well as being able to be used as a scripting language [12].

## Hello World

**Java:**

```
// Hello World in Java
class HelloWorld {
    static public void main(String args[]) {
        System.out.println("Hello World!");
    }
}
```

**C++:**

```
// Hello World in C++
#include <iostream.h>
Main() {
    cout << "Hello World!" << endl;
    return 0;
}
```

**Python:**

```
# Hello World in Python
print("Hello World!")
```

Figure 4.1: A comparison between Java, Python and C++ to print an output to the console. [1]

Python gets considered to be easy to learn the language due to its high readability and is recommended by many exam boards as the language to use for teaching Computer Science at GCSE and A-Level level [?]. Python's simple and easy to learn syntax emphasises on readability, which, as a result, reduces the cost of program maintenance [12, 13].

Python gets compared to a lot of other languages. However, due to the requirements and expectations of the application, we will compare it to other similar style applications that can potentially do a similar job. These being Java, JavaScript and C++. In general, the choice of the programming language to use is many other real-world constraints, for example, financial cost, availability, training and even personal preferences and attachments. However, we will focus on language issues for the comparisons.

In comparison to Java, Python programs will typically take 3-5 times quicker (See fig: 4.1) to develop but will have a slower run time. The time difference gets attributed to Python's built-in data types and its dynamic typing [14]. As Java gets better characterised as a low-level implementation language, this would be the language of choice if application execution speed was the deciding factor. If this is not a factor, then there is no real benefit over Python.

What gets said about Java is also the same when comparing C++ to Python. It is often



5-10 times shorter than equivalent C++ code. Anecdotal evidence suggests that one Python programmer can finish in two months what two C++ programmers cannot complete in a year. Python shines as a glue language, used to combine components written in C++ [14].

In comparison to JavaScript (JS), Python's 'object-based' subset is very similar to JS. Python supports a programming method that uses simple variables and functions, similar to JS, that do not need class definitions. However, Python also supports writing for much larger programs, which leads to more reusable code by using an accurate OOP way while with JS, that is all that it can do [14].

Another language that presented itself to us was the R language. R is a language and environment for statistical computing and graphics. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues [15]. Many users think of R as a statistics system [15]. Academics and statisticians have developed R over two decades. There are around 12000 packages available in CRAN (open-source repository). The wide variety of library makes R the first choice for statistical analysis, especially for specialised analytical work [16].

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering) and graphical techniques, and is highly extensible. One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed [15]. The cutting-edge difference between R and the other statistical products is the output. R has fantastic tools to communicate the results. Rstudio comes with the library knitr. Communicating the findings with a presentation or a document is easy [16].

R and Python are both open-source programming languages with a large community. New libraries or tools are added continuously to their respective catalogue. R is mainly used for statistical analysis, while Python provides a more general approach to data science. R and Python are both state of the art in terms of programming language oriented towards data science. Learning both of them is, of course, the ideal solution. R and Python requires a time-investment, and such luxury is not available for everyone. Python is a general-purpose language with a readable syntax. R, however, is built by statisticians and encompasses their specific language [16].

Python can pretty much make the same tasks as R: data wrangling, engineering, feature selection web scrapping, creating an app, for example. Python is a tool to deploy and implement

#### 4. *Implementation*

---

machine learning at a large-scale. Python codes are easier to maintain and more robust than R. Years ago; Python did not have many data analysis and machine learning libraries.

Recently, Python is catching up and provides cutting-edge API for machine learning or Artificial Intelligence. Most of the data science job can get done with five Python libraries: Numpy, Pandas, Scipy, Scikit-learn and Seaborn [16].

Python, on the other hand, makes replicability and accessibility easier than R., if we need to use the results of our analysis in an application or website, Python is the best choice [16].

In 2019 there was an active number of 26.66 billion devices attached to the internet [17,18], with an estimation of 35 billion in 2021 [17] and by 2025 75.44 billion [18]. Experts estimate that the IoT device market will reach \$1.1 trillion in 2026 [17]. Every Second 127 new devices get connected to the world wide web [17].

With so many devices on the internet, an important consideration we had was to make the application web-based. By creating the application for the internet, this would allow potentially many more people to be able to access the application and interact with the different ML models.

JS gets regarded as more of the language of the world-wide-web [19]. It got initially designed to be used client-side in a web browser. However, it has in more recent years started to branch out and be able to be used to create applications on, not only the front end of the web but also desktops, servers and mobile platforms natively. For example, React Native, Node.js and TypeScript. JavaScript is also incredibly useful, allowing developers to be able to create apps with audiences in the millions quickly [20].

The decision on what language to use we a close call between Python and JavaScript, this was due to the massive amounts of libraries that were on offer and the support communities that were in place. With both being open source and both having essential libraries available to interact with machine learning models and visualisation tools, both could have been a perfect fit for the intended application. However, we decided upon using Python. Python was chosen based on it being the go-to language for anything machine learning related, and its ability to be able to be used multiplatform on desktops or mobile devices. Python supports modules and packages, which encourages programs to be developed modularity and therefore allows code to get reused. The Python interpreter and the extensive standard library are available in source or binary form without charge for all major platforms and can get freely distributed [12]. There was also an additional factor that the author was more familiar with Python and its required libraries compared to the libraries that will get required for using JavaScript.

There was the additional decision to use HTML and CSS within a small part of the project, the 'Learning Zone', based on the quickness of being able to create and host the webpages containing the learning content. It was allowing the learning content to evolve without having any impact on the overall development of the main application, allowing the learning content to be an individual entity within the main application.

## 4.2 Frameworks

### 4.2.1 GUI Framework

With the nature of the application, we needed to make the application have a Graphical User Interface (GUI). Having a GUI allowed the learning to be a lot more hands-on and allow the players to see what is happening within the models, especially when they interact with them.

Therefore, due to the GUI requirement, three GUI libraries presented themselves to us. These were Pygame, PyQt5 and Tkinter.

Pygame is a free, open-sourced library. Released under the LGPL licence, Pygame is a set of Python modules designed for writing video games. Pygame adds functionality on top of the standard Python library. Pygame allows the user to create fully featured games and multimedia programs in the python language [21].

Pygame is highly portable and runs on nearly every platform and operating system, and it gets downloaded millions of times [21].

With the main aim of the application to be a game, Pygame was a strong contender. It was providing modules that can handle a lot of the key gaming mechanics and multiple screen switching. However, it lacked some key features that were deemed essential for the application. It was unable to provide a library that could create interactable graphs to be used as data inputs for the models and be able to render HTML and CSS content for the Learning Zone. Therefore reducing the amount of flexibility, it got decided upon for using HTML and CSS for the learning content. Therefore, meaning that all the content would need to be hardcoded. If any changes were needed, a significant transformation would need to happen to the overall code, instead of just changing the web content.

PyQt is a set of Python v2 and v3 bindings for The Qt Company's Qt application framework and runs on all platforms supported by Qt including Windows, macOS, Linux, iOS and Android. PyQt5 supports Qt v5. PyQt4 supports Qt v4 and will build against Qt v5. The bindings are implemented as a set of Python modules and contain over 1,000 classes [13].

PyQt brings together the Qt C++ cross-platform application framework and the cross-platform interpreted language Python. Qt is more than a GUI toolkit. It includes abstractions of network sockets, threads, Unicode, regular expressions, SQL databases, SVG, OpenGL, XML, a fully functional web browser, a help system, a multimedia framework, as well as a rich collection of GUI widgets. Qt classes employ a signal/slot mechanism for communicating between objects that is type safe but loosely coupled making it easy to create re-usable software components [13].

Qt also includes Qt Designer, a graphical user interface designer. PyQt is able to generate Python code from Qt Designer. It is also possible to add new GUI controls written in Python to Qt Designer [13].

PyQt combines all the advantages of Qt and Python. A programmer has all the power of Qt but can exploit it with the simplicity of Python [13].

Tkinter is the third option. Tkinter commonly comes bundled with Python, using Tk and is Python's standard GUI framework. It is famous for its simplicity and graphical user interface. It is open-source and available under the Python License [22].

Tkinter is Python's de-facto standard GUI (Graphical User Interface) package. It is a thin object-oriented layer on top of Tcl/Tk. Tkinter is not the only GuiProgramming toolkit for Python. It is however the most commonly used one. CameronLaird calls the yearly decision to keep TkInter "one of the minor traditions of the Python world [23]."

Tkinter supports functionality with Matplotlib, with Matplotlib offering libraries to allow handling the backend of the graph creation interacting with the GUI library. However, unlike QT, Tkinter does not support any GUI designer. Therefore the GUIs will have to be created programmatically, which will give more control, might involve more of a learning curve and potentially more time to implement in the initial stages.

After reviewing the different GUI libraries, PyQt was the decided library to use. We believed it would give us the ability to have

### 4.3 Packages

### 4.4 IDE

PyCharm is a dedicated Python Integrated Development Environment (IDE) providing a wide range of essential tools for Python developers, tightly integrated to create a convenient environment for productive Python, web, and data science development [24]. While PyCharm is

a very popular IDE, and one that we have had experience with before, it is not, however, one that we have had many experiences using compared to other IDEs. While it does provide much functionality and it is a lot easier to use and keep our directories organised compared to Python's provide IDE, it has, however, not been an IDE that has flowed well when we have used it.

Visual Studio Code (VS Code) is a free source-code editor made by Microsoft for Windows, Linux and macOS. Features include support for debugging, syntax highlighting, intelligent code completion, snippets, code refactoring, and embedded Git. Visual Studio Code combines the simplicity of a source code editor with powerful developer tooling, like IntelliSense code completion and debugging. First and foremost, it is an editor that gets out of the user's way. The delightfully frictionless edit-build-debug cycle means less time fiddling with the required environment, and more time executing ideas [25].

Microsoft claims that VS Code, at its heart, lightning-fast code features a lightning-fast source code editor, which is perfect for day-to-day use. With support for hundreds of languages, VS Code helps the user be instantly productive with syntax highlighting, bracket-matching, auto-indentation, box-selection, snippets, and more [25]. For serious coding, the user will often benefit from tools with more code understanding than just blocks of text. Visual Studio Code includes built-in support for IntelliSense code completion, rich semantic code understanding and navigation, and code refactoring [25]. Which we can say from experience is mostly true. However, on occasions, it has provided code completion that was not intended or needed. VS Code also allows the user to customise every feature to their liking and install any number of third-party extensions. While most scenarios work "out of the box" with no configuration, VS Code also grows with you [25]. Which, from our experience, we can say is true. VS Code has grown with us. The VS Code community has provided many extensions that have helped with our workflow.

Atom is developed and released by GitHub [26]. Atom is free, and an open-sourced code editor. Atom is a self-labelled 'a hackable text editor for the 21st century'. Atom, like VS Code, allows developers to fully customise the look, feel, and requirements to speed up their workflows.

However, Atom still allows developers to use it productively without ever touching a config file. Atom comes pre-loaded with eight syntax themes and four UI, two light and two dark, but if none of them provides any interest, Atom makes it easy and quick to install customised themes created by a third-party or to create one [26]. However, apart from pre-created extensions to help with code linting and code autocomplete abilities, none of these features is of any

interest to us. The main factor does the IDE have a friendly UI and does it seem not to hinder our workflow. Which is safe to say, it does have a friendly UI and does not hinder our workflow at all.

After trailing the different IDEs, we believed that the best option going forward was the VS Code IDE. We have chosen this IDE because of two key factors. The first one being that it supported all the libraries needed, whether it was pre-installed or through downloading additional extensions, and that we have had a better familiarity with the IDE's interface from previous uses and projects.

## 4.5 Intricacies of the Game Components

### 4.5.1 Gameplay Area

[Need to Finish in application]

### 4.5.2 Learning Zone Area

The Learning Zone (LZ) area is an area that we intend to allow the user to do most of the learning. The LZ, in terms of UI, is very basic. It has a web browser window and three buttons. The web browser window is where the HTML and CSS documents, which the created web documents, get displayed within the application.

The web document consists of a welcome page, outlining the content, a "What is Machine Learning?", "Task Driven vs Data-Driven", "Supervised and Unsupervised Learning", "Classification", "Support Vector Machines (SVM)", "k-Nearest Neighbour", "Neural Networks", "Regression", "Linear Regression", "Logistic Regression", "Clustering", "K-Means", "Gaussian Mixture Model", "Dimensionality Reduction", "Principal Component Analysis", "Linear Discriminant Analysis" and "Association Rule" web pages.

The web pages follow a similar layout design. A blue background, a yellow background layer on top with an offset grey colour behind the text. Each page contains title at the top with a dark grey background. The content of the web pages either we an overview, for example, "Clustering", which looked into clustering as a whole and what the different types were. Alternatively, a web page would explain a specific algorithm, for example, "K-Means", which explained the intricacies of how the algorithm worked and the critical mechanics behind it.

The three buttons at the bottom of the application screen trigger three different actions. All of which match to the intended buttons, a home button, to go back to the main menu, a free

play button to send the player to the free play area with the intended algorithm that the user was learning about, and a quiz button which loads a multi-choice quiz.

When the player clicks the free play button, the application checks the HTML documents title tag and loads up the required model in the free play section. However, if the player clicks the button and the web page does not have a model available, within the free play section or it is just a general overview page, a message box will appear. The message box intension is to let the user know that they can not progress to the free play zone and a list of the available models (see fig: ??).

The Quiz area is an additional area to the learning zone. The Quiz area is where the user can get tested on what they learned in the LZ area, in the form of a multiple-choice quiz. When the user is viewing a topic on the LZ, and they decide to take a quiz, the user will click on the quiz button, and this will read the title tags of the HTML document and open the required quiz. The quiz questions and answers are within their text file, and the name of the file matches the title tag's content. The text file itself holds the information in the format of a 2D array, that has the question at position zero, the answer at position one and then position 2 to 5 are the multiple-choice options. This information from the text file populates a question label and four buttons, allowing the user to click what button they think is the answer. The total of correct answers get added up and displayed to the user at the end in a message box.

### **4.5.3 Free Play Area**

The Free Play (FP) zone is an area where the user gets to interact and play with different ML models. The models include Linear Regression, K-Means, Neural Networks, Linear Discriminant Analysis, Gaussian Mixture Models and SVM. The intension for the FP zone was to have all the models explained in the learning zone be available for the player to interact with, so it could help them fully understand how the model works by allowing the user to manipulate parameters and data points. However, due to time restrictions, there are only six models available, with 5 of the models having real interactivity but to different degrees.

When the FP zone is accessed, unless accessed through the Learning Zone, a randomly selected model gets displayed to the user from the list mentioned before. On first glance, the user has multiple areas to either interact with or present information to them. The screen has a Widget that is linked to a Matplotlib library to handle PyQt5 backend interactions. Also, a model overview is displayed next to the widget, it tells the user information about the model, for example, the type of learning it is, supervised or unsupervised, the name of the model and a

#### 4. Implementation

---

brief overview of the model. Just beneath the widget and overview is a group box that contains all the settings for the model and data interaction. The model settings group box contains combo boxes, radio buttons, checkboxes, line edits and buttons, which all do different things depending on the model and data sets selected. Within the model settings group box, there are three additional group boxes. These are 'Model Attribute(s)', 'Model Parameter(s)' and 'Data Options' with each group box displaying different content depending upon the model and data options selected in the combo boxes.

The model combo box contains six values, and these are 'Please Select', 'K-Means', 'LDA', 'Linear Regression', 'GMM', 'SVM', 'Neural Networks'. Once one of these options is selected, apart from the 'Please Select', the desired model will display in the Matplotlib widget area. The Model attributes and parameters boxes will display the required information unless the models 'LDA', 'SVM' and 'GMM' are selected. Instead, a label placeholder saying, 'No Options available, yet!' will be displayed. While LDA has a fully interactive model in the Matplotlib widget, it does not present options for the user to change within the model, the user can only click on the widget and place points, which the model will then apply and create the required actions. Therefore a place holder label appears stating to click in the game widget to interact with the model. While LDA and GMM both have the ability for the model visualisations to toggle on and off, showing how the models have fit their data, GMM has little much additional functionality. GMM only allows the user to toggle on and off the visualisation, which is the model predicting the 'Iris' dataset clusters. However, SVM only displays the model's output, again using the 'Iris' data set, but the output shows the boundary lines and area that each partition covers.

While on the other hand, the Linear Regression, K-Means and Neural Network models display different options. Linear Regression displays to the user labels in the attributes group box to show them the values for the intercept, estimated coefficient and outcome. There is also a line edit available for the user to input a value and see what the model would predict out, which gets displayed in the output label. However, Linear Regression does not have any model parameters, and this is due to the values getting deemed as not having much impact on the model and limited implementational time.

When K-Means gets selected, both the attributes and parameters group boxes have information and selection options displayed to the user. The attributes group box displays the information for Inertia, the number of iterations that got performed fitting the data, prediction, which relies on the user to click within the Matplotlib widget and the X and y coordinates get



displayed along with a cluster prediction label in the output label. There is also a distance from the centroid value displayed, and this value got achieved by using the SKLearn Metrics library. The model parameters group box displays multiple line edits and a combo box that allows the user to input values to the model. These will alter the K-Means k value (number of clusters), the number of initialisers, the max number of iterations and the underlying algorithm (auto, full or Elkan), that gets used. The k value is independent of the number of clusters in the data options, so they do not impact on each other. Allowing the k value to be changed independently will allow the user to be able to experiment with the model to see how two, three or other k values affect the prediction, even when known that the data may have, for example, five different clusters. K-Means also brings up an additional option, and this is to be able to switch on and off the centres of the clusters. When the checkbox gets enabled, this will lay on top of the data points an 'X' where each of the cluster centres is and when it is disabled, it will remove the 'X'.

[NN Att no. of layers/ neurons -> params set the values. Not implemented in-app yet!]

The data combo box displays the data options for the different models and depending on what option is selected depends on the information that is on offer to the user in the data options group box. If the custom data option is selected, then the group box will display different options to the user to be able to generate custom data points to be displayed in the Matplotlib widget game screen. The only models that have this option are the Linear Regression and K-Means models. Linear Regression has the data options 'Diabetes' and 'Boston House Prices' and K-Means has the data options 'Iris' and 'Moons'. When these are selected, the data option displays radio button options for the user to select the features they would like the model to display, on the X and y-axis, and get fitted. Linear Regression's 'Custom' data option allows the user the ability to change the random generated data's settings. These settings include the number of data samples and if there are any outliers wanted, if so then an option to add the number of outliers. Whereas K-Means 'Custom' data option allows the user the ability to change the number of clusters to generate and the number of data samples wanted. The other models have a label placeholder saying, 'no data options selected yet!' and no actual data selection options. In the case of LDA and Neural Network, this is more due to the decision we made. Based on the way the model's fit function was implemented, not allowing the user to generate random data was decided. Doing so would impact on how the model gets interacted with by the user. However, in the case of the other models, it was a lack of time that impacted the inability to add this feature. Though, it was always the intention to add it.

#### *4. Implementation*

---

The final aspects of the Model options group box are three buttons, a play button, a clear button and a home button. Where the home button is self-explanatory in terms of returning the user to the main menu, and the clear button resetting the Matplotlib Widget axis contents. The play button is where the primary handling of how the user interacts with the back end of the models and datasets. When the user inputs information, they are required to press the play button for these features to be implemented.

##### **4.5.4 Achievements Area**

The intensions for the Achievement Area was displaying to the user all of the gamification badges available. Providing an overview and hints on how to unlock them. The achievements were going to have a bronze, silver and gold level, and we planned to be unlocked once the user had completed specific tasks like playing the game three times or completing a quiz. However, due to time limitation, this was not possible, and the application, as it currently stands, displays a coming soon image and button for the user to navigate back to the main menu.

#### **4.6 Example user stories (A UML term for case studies or example playthroughs)**

## Chapter 5

# Conclusions and Future Work

In this document we have demonstrated the use of a LaTeX thesis template which can produce a professional looking academic document.

### 5.1 Contributions

The main contributions of this work can be summarized as follows:

- **A LaTeX thesis template**

Modify this document by adding additional top level content chapters. These descriptions should take a more retrospective tone as you include summary of performance or viability.

- **A typesetting guide of useful primitive elements**

Use the building blocks within this template to typeset each part of your document. Aim to use simple and reusable elements to keep your document neat and consistently styled throughout.

- **A review of how to find and cite external resources**

We review techniques and resources for finding and properly citing resources from the prior academic literature and from online resources.

### 5.2 Future Work

Future editions of this template may include additional references to Futurama.



# Bibliography

- [1] Things Tech, “Which programming language to start with as a beginner,” 2020, [Online; accessed August 10th, 2020]. [Online]. Available: <https://thingsteck.wordpress.com/>
- [2] Internet Live Stats. (2020). [Online]. Available: <https://www.internetlivestats.com>
- [3] G. Branwen. (2020) Internet search tips. [Online]. Available: <https://www.gwern.net/Search>
- [4] RELX Group. (2019) Mendeley. [Online]. Available: <https://www.mendeley.com>
- [5] J. T. Kajiya, “The rendering equation,” in *Proceedings of the 13th Annual Conference on Computer Graphics and Interactive Techniques*, ser. SIGGRAPH ’86. New York, NY, USA: ACM, 1986, pp. 143–150.
- [6] Apple Inc. (2020) Swift. [Online]. Available: <https://developer.apple.com/documentation/swift>
- [7] M. Potuck. (2020) Apple hits 1.5 billion active devices with 80iphones and ipads running ios 13. [Online]. Available: <https://9to5mac.com/2020/01/28/apple-hits-1-5-billion-active-devices-with-80-of-recent-iphones-and-ipads-running-ios-13/>
- [8] K. Finley. (2020) Python is more popular than ever. [Online]. Available: <https://www.wired.com/story/python-language-more-popular-than-ever/>
- [9] B. Popper. (2020) The 2020 developer survey results are here! [Online]. Available: <https://stackoverflow.blog/2020/05/27/2020-stack-overflow-developer-survey-results/>
- [10] A. Goel. (2020) Best programming language to learn in 2020 (for job and future). [Online]. Available: <https://hackr.io/blog/best-programming-languages-to-learn-2020-jobs-future>

- [11] Stack Overflow. (2013) Python vs cpython. [Online]. Available: <https://stackoverflow.com/questions/17130975/python-vs-cpython>
- [12] Python Software Foundation. (2020) What is python? executive summary. [Online]. Available: <https://www.python.org/doc/essays/blurb/>
- [13] Riverbank Computing. (2020) What is pyqt? [Online]. Available: <https://riverbankcomputing.com/software/pyqt/intro>
- [14] Python Software Foundation. (2020) Comparing python to other languages. [Online]. Available: <https://www.python.org/doc/essays/comparisons/>
- [15] The R Foundation. (2020) What is r? [Online]. Available: <https://www.r-project.org/about.html>
- [16] Guru 99. (2020) R vs python: What's the difference? [Online]. Available: <https://www.guru99.com/r-vs-python.html#:~:text=New%20libraries%20or%20tools%20are,general%20approach%20to%20data%20science.&text=Python%20is%20a%20general%2Dpurpose,and%20encompasses%20their%20specific%20language.>
- [17] G. D. Maayan. (2020) The iot rundown for 2020: Stats, risks, and solutions. [Online]. Available: <https://securitytoday.com/Articles/2020/01/13/The-IoT-Rundown-for-2020.aspx?Page=1>
- [18] Statista. (2020) Internet of things (iot) connected devices installed base worldwide from 2015 to 2025. [Online]. Available: <https://www.statista.com/statistics/471264/iot-number-of-connected-devices-worldwide/>
- [19] World Wide Web Foundation. (2020) History of the web. [Online]. Available: <https://webfoundation.org/about/vision/history-of-the-web/>
- [20] T. DeGroat. (2019) The history of javascript: Everything you need to know. [Online]. Available: <https://www.springboard.com/blog/history-of-javascript/>
- [21] Pygame. (2020) About. [Online]. Available: <https://www.pygame.org/wiki/about>
- [22] A. Sharma. (2020) Introduction to gui with tkinter in python. [Online]. Available: [https://www.datacamp.com/community/tutorials/gui-tkinter-python?utm\\_source=adwords\\_ppc&utm\\_campaignid=](https://www.datacamp.com/community/tutorials/gui-tkinter-python?utm_source=adwords_ppc&utm_campaignid=)

898687156&utm\_adgroupid=48947256715&utm\_device=c&utm\_keyword=&utm\_matchtype=b&utm\_network=g&utm\_adposition=&utm\_creative=229765585186&utm\_targetid=aud-299261629574:dsa-429603003980&utm\_loc\_interest\_ms=&utm\_loc\_physical\_ms=1007460&gclid=Cj0KCQjwsuP5BRCoARIsAPtX\_wGxaCEuqKDVDDyVASqiCw2zIKYwN2Duo3DObWGcfwQAjH3oxWK5WfoaAhafEALw\_wcB

- [23] Python Software Foundation. (2020) Tkinter. [Online]. Available: <https://wiki.python.org/moin/TkInter>
- [24] PyCharm. (2020) Get started. [Online]. Available: <https://www.jetbrains.com/help/pycharm/quick-start-guide.html>
- [25] Microsoft. (2020) Why did we build visual studio code? [Online]. Available: <https://code.visualstudio.com/docs/editor/whyvscode>
- [26] CloudApp. (2020) A guide to atom text editor. [Online]. Available: <https://www.getcloudapp.com/blog/how-to-use-atom-text-editor>





## Appendix A

# Implementation of a Relevant Algorithm

```
1 #include <stdio.h>
2
3 int main(int argc, char *argv[]) {
4     printf("Hello world.\n");
5     return 0;
6 }
```

Listing A.1: An implementation of an important algorithm from our work.



## **Appendix B**

# **Supplementary Data**

The results of large ablation studies can often take up a lot of space, even with neat visualization and formatting. Consider putting full results in an appendix chapter and showing excerpts of interesting results in your chapters with detailed analysis. You can use labels and references to refer the reader here for the full data.